

Predicting diverse M-best protein contact maps

Siqi Sun, Jianzhu Ma, Sheng Wang and Jinbo Xu

Toyota Technological Institute at Chicago,

Chicago, IL, USA, 60637

Email: siqi.sun, majianzhu, wangsheng@ttic.edu, jinbo.xu@gmail.com

Abstract—Protein contacts contain important information for protein structure and functional study, but contact prediction from sequence information remains very challenging. Recently evolutionary coupling (EC) analysis, which predicts contacts by detecting co-evolved residues (or columns) in a multiple sequence alignment (MSA), has made good progress due to better statistical assessment techniques and high-throughput sequencing. Existing EC analysis methods predict only a single contact map for a given protein, which may have low accuracy especially when the protein under prediction does not have a large number of sequence homologs. Analogous to *ab initio* folding that usually predicts a few possible 3D models for a given protein sequence, this paper presents a novel structure learning method that can predict a set of diverse contact maps for a given protein sequence, in which the best solution usually has much better accuracy than the first one. Our experimental tests show that for many test proteins, the best out of 5 solutions generated by our method has accuracy at least 0.1 better than the first one when the top $L/5$ or $L/10$ (L is the sequence length) predicted long-range contacts are evaluated, especially for protein families with a small number of sequence homologs. Our best solutions also have better quality than those generated by the two popular EC methods Evfold and PSICOV.

I. INTRODUCTION

Protein contacts contain important information for protein folding and recent works indicate that one correct long-range contact for every 12 residues in the protein allows accurate topology level modeling [20]. Thanks to high-throughput sequencing and better statistical assessment techniques, evolutionary coupling (EC) analysis for contact prediction has made good progress, which makes *de novo* structure prediction of some large proteins possible [25], [28], [34]. Nevertheless, contact prediction accuracy for many proteins is still low even if only the top $L/10$ (L is the sequence length) predicted contacts are evaluated.

A number of contact prediction methods have been developed including [6], PSICOV [10], Evfold [25], SVMSEQ [41], NNcon [35], SVMcon [7], CMAPpro [11], PhyCMAP [39] and Astro-Fold [21]. Some of them are pure EC-based unsupervised methods while others are supervised machine learning methods that integrate both EC and a variety of non-EC information. This paper focuses on the pure EC-based unsupervised approach, exploring one more way to improve residue evolutionary coupling (EC) analysis. EC analysis detects co-evolved residues from the MSA (multiple sequence alignment) of a protein family and then predicts

two co-evolved residues to form a contact. This approach is based upon an observation that a pair of co-evolved residues is often found to be spatially close in the 3D structure. Along with many more sequences generated by high-throughput sequencing, some sophisticated global statistical methods (or structure learning methods), such as maximum entropy and probabilistic graphical models, are developed to infer residue co-evolution from MSA [3], [8], [23], [25], [40]. These global statistical methods can differentiate direct from indirect residue couplings and thus, are more accurate than the traditional mutual information, which is a local statistical method, in predicting contacts. Global statistical methods differ from local methods in that the former considers their relationship with all the other residues in determining if two residues are co-evolved or not while the latter ignores this kind of relationship. Representative tools of recent global statistical methods for contact prediction include Evfold [25], PSICOV [18] and GREMLIN [3]. See [9] for an excellent review of EC analysis.

Although quite a few new contact prediction methods have been developed recently, existing methods predict only a single contact map (also called contact graph) for a protein sequence. Such a single solution may have low accuracy especially when the protein under prediction does not have a large number of sequence homologs, i.e., no sufficient information for the determination of the ground-truth contact map. This paper tackle the challenging contact prediction problem from another perspective. Analogous to *ab initio* folding that usually predicts a few possible 3D models for a given protein sequence, this paper presents a novel structure learning method that can predict M (a small constant) diverse contact maps for a given protein sequence such that the best of the predicted contact maps usually has much better accuracy than the first one. In addition to *ab initio* folding, the idea of generating multiple solutions has been studied for other bioinformatics problems such as sequence and structure alignment and genome assembly, but not for contact prediction. All existing contact prediction methods generate only a single contact map for a protein sequence, although they may assign a confidence score to each predicted contact.

Our experimental tests show that for many test proteins, we only need to generate five different contact maps for a single protein in order to yield a contact map with much

better accuracy than the first one and those generated by two popular EC analysis methods PSICOV and Evfold. This is especially true when the protein family under prediction has a small number of sequence homologs. This paper focuses on how to generate diverse M -best contact maps, but does not study how to re-rank them. A possible way to do so is that we may conduct protein folding simulation using each predicted contact map and then select the one resulting in the best protein 3D models according to some atom-level energy functions or model quality assessment methods.

II. METHODS

A. Background

In this section we introduce some notations and PSICOV [18], one state-of-the-art co-evolution analysis method for contact prediction. PSICOV formulates the contact prediction problem using Gaussian Graphical model (GGM). Our proposed method will also be based upon GGM, but the idea we propose here can be easily adapted to other models, such as the Potts model [14] and the pseudo-likelihood model [31].

We run the buildali.pl program in the HHpred package [33] to find its sequence homologs and then build an MSA of these homologs. We denote the MSA as $X = (x_1, x_2, \dots, x_N)$, where N is the number of sequence homologs and each sequence $x_i = (x_{i1}, x_{i2}, \dots, x_{iL})$ is a vector with length L . Each x_{ij} takes a value from a set Φ with cardinality 21 consisting of 20 amino acids and one gap symbol. So, x_{ij} is a vector of 21 binary variables and each element in x_{ij} indicates the presence or absence of a specific amino acid at row i and column j .

Let $f_i(a)$ denote the frequency of amino acid a at position i and $f_{ij}(a, b)$ the frequency of one pair of amino acids a and b at positions i and j . We may calculate the $21L \times 21L$ covariance matrix for the MSA as follows.

$$S_{ij}^{ab} = f_{ij}(ab) - f_i(a)f_j(b), \quad (1)$$

where the value of S_{ij}^{ab} implies the correlation of amino acids a and b at positions i and j . In this paper we also use S_{ij} to denote the 21×21 sub-matrix corresponding to two positions (columns) i and j , and the value of the submatrix at position (a, b) is S_{ij}^{ab} .

Given the empirical (or data) covariance matrix S above, we can compute its inverse by minimizing the negative log-likelihood with a l_1 penalty term:

$$\begin{aligned} \Omega &= \arg \min_W F(W) + \lambda \|W\| \\ &= \arg \min_W \text{tr}(SW) - \ln \det(W) + \lambda \sum_{ij} |W_{ij}| \end{aligned} \quad (2)$$

where $F = \text{tr}(SW) - \ln \det(W)$ is the negative log-likelihood of the observed sequences in an MSA, Ω is the estimated precision matrix (i.e., inverse of S), and λ is a hyper-parameter controlling the sparsity level of Ω .

The structure of the resultant precision matrix Ω encodes the residue interaction pattern and can be used to derive a contact map [18]. To efficiently optimize Eq.(2), please refer to [16].

PSICOV formulates the contact prediction problem using Eq. (2) and predicts contacts from the estimated precision matrix using the l_1 norm of the submatrix corresponding to two positions (or columns) i and j (excluding the contribution of gaps), i.e.

$$\tilde{\Omega}_{ij} = \sum_{a,b \in \Phi} |\Omega_{ij}^{ab}|, \quad (3)$$

where a and b range over the 20 amino acids. That is, $\tilde{\Omega}_{ij}$ is a measure of interaction strength between positions i and j . The higher $\tilde{\Omega}_{ij}$, the more likely there is a contact between i and j . PSICOV also reduces the entropic and phylogenetic bias by correcting $\tilde{\Omega}_{ij}$ with average product correction (APC) [13]. Therefore the final prediction score can be represented as

$$\tilde{\Omega}_{ij}^{\text{corrected}} = \tilde{\Omega}_{ij} - \frac{\tilde{\Omega}_i \cdot \tilde{\Omega}_j}{\tilde{\Omega}_{..}} \quad (4)$$

where $\tilde{\Omega}_i$ is the mean of all the $\tilde{\Omega}_{ij}$ sharing position i , $\tilde{\Omega}_j$ is the mean of all the $\tilde{\Omega}_{ij}$ sharing position j and $\tilde{\Omega}_{..}$ is the mean of all the $\tilde{\Omega}_{ij}$. Finally, those residue pairs with large $\tilde{\Omega}_{ij}^{\text{corrected}}$ are predicted to form a contact.

B. A new model for predicting diverse M -best contact maps

For many proteins, a single solution may deviate a lot from the ground truth [5]. Here we propose a method that can generate a set of alternative contact maps for a protein under prediction. We require that the set of solutions to be generated are diverse, so that the best solution may have much better accuracy. Note that since the search space is huge, it is still challenging to yield a solution with better accuracy even if we predict a small number (e.g., 5 or 10) of diverse solutions.

We generate a set of diverse solutions using an iterative process. Let Ω^1 denote the first solution to Eq.(2), i.e., the precision matrix generated by the popular method PSICOV. Then we generate a new solution Ω^2 based on Ω^1 , such that the distance between them, $d(\Omega^1, \Omega)$, is larger than a desired value ϵ , which quantifies the diversity of the alternative solutions. Next we generate Ω^3 based on both Ω^1 and Ω^2 . Generally speaking, in step $m+1$ we generate Ω^{m+1} based on all previous precision matrices $\Omega^1, \Omega^2, \dots, \Omega^m$, such that the new solution is different from all of them. More specifically, we will generate a new solution by solving the following optimization problem:

$$\begin{aligned} \min_{\Omega} F(\Omega) + \lambda \sum_{ij} |\Omega_{ij}| \\ \text{s.t. } d(\Omega^k, \Omega) \geq \epsilon, k = 1, \dots, m, \end{aligned} \quad (5)$$

where $d(\Omega^k, \Omega)$ measures the difference between two alternative solutions Ω^k and Ω . The new solution resulting from this formulation minimizes the negative log-likelihood of the observed sequences in the MSA, subject to the constraint that it must be different from previous solutions by a given cutoff ϵ . Note that even though the previous results $\Omega_1, \Omega_2, \dots, \Omega_m$ are not explicitly expressed in the objective function, they impact the new solution Ω through the constraints.

We use the popular Hamming distance to measure the difference between two solutions, which is defined as follows.

$$d(\Omega^k, \Omega) = (L^2 - \langle \delta_0(\Omega^k), \delta_0(\Omega) \rangle) / 2 \quad (6)$$

where L^2 is the number of residue pairs (or submatrices) in Ω , and δ_0 maps a real-valued $21L \times 21L$ matrix to a binary $L \times L$ matrix. More precisely, if the submatrix Ω_{ij} for residue pair (i, j) is not 0, then $\delta_0(\Omega_{ij})$ takes value 1; otherwise -1. The operator $\langle \cdot, \cdot \rangle$ calculates the inner product of two matrices. The Hamming distance is always larger than or equal to zero. The smaller the distance is, the more similar contact maps implied by the two solutions.

However, it is challenging to optimize Eq.(5) with respect to this distance constraint function since $\delta_0(\Omega)$ is neither convex nor differentiable. To tackle this problem, we relax the Hamming distance as follows.

$$\hat{d}(\Omega^k, \Omega) = -\langle \delta_0(\Omega^k), |\Omega| \rangle = -\sum_{ij} \delta_0(\Omega^k)_{ij} |\Omega_{ij}|, \quad (7)$$

Since at this point Ω^k is a constant, the above relaxation is actually a negative re-weighted l_1 norm of the precision matrix Ω . Such a relaxation has the following desired properties:

- If there is a contact between i and j in the previous solution Ω^k , the chance of there being a contact between i and j in Ω will be decreased.
- If there is no contact between i and j in all previous Ω^k , the chance of there being a contact between i and j in Ω will be increased.
- The optimization is also easier since the relaxation is a convex function.

The first two properties ensure that the new solution will be different from all previous ones. The third property will not only simplify the optimization, but also significantly reduce the computational cost by using a specific procedure to be described in the next section. Finally, our original formulation (5) becomes

$$\min_{\Omega} F(\Omega) + \lambda \sum_{ij} |\Omega_{ij}| \quad (8)$$

$$s.t. \quad \hat{d}(\Omega^1, \Omega) \geq \epsilon, k = 1, \dots, m$$

Where the continuous distance function \hat{d} replaces the 0-1 Hamming distance d in (5).

C. Optimization

It is easy to prove that for every ϵ (details ignored due to space limit), there exists some $\mu(\epsilon)$ such that the problem (8) is equivalent to minimizing

$$G(\Omega) = F(\Omega) + \lambda \|\Omega\|_1 + \sum_{k=1}^m -\mu_k(\epsilon) \hat{d}(\Omega^k, \Omega) \quad (9)$$

$$= tr(S\Omega) - \ln \det(\Omega) + \sum_{ij} (\lambda - \sum_{k=1}^m \mu_k(\epsilon) \delta_0(\Omega^k)_{ij}) |\Omega_{ij}|$$

Compared to (8), the new formulation (9) is appealing since it has no constraints and also has the same form as the original formulation (2), which has been extensively studied and many solvers are developed. This allows us to apply existing solvers to solve our new formulation (9). The hyper parameters $\mu_k(\epsilon)$ controls how far away we want our solution to deviate from the previous solution Ω^k . To prevent overfitting, we set all μ_k to be the same.

Note that all the μ_k shall be upper bounded by λ/m . Otherwise, $\lambda'_{ij} = \lambda - \sum_{k=1}^m \mu_k \delta_0(\Omega^k)_{ij}$ becomes negative and thus, formulation(9) becomes non-convex. The intuition underlying this upper bound is that the difference between two solutions shall not be too big, which makes sense since all the solutions shall minimize the negative log-likelihood of the observed sequences in an MSA.

Finally, we generate a new contact map by minimizing the following function.

$$G(\Omega) = tr(S\Omega) - \ln \det(\Omega) + \sum_{ij} (\lambda - \sum_{k=1}^m \mu \delta_0(\Omega^k)_{ij}) |\Omega_{ij}|$$

$$= tr(S\Omega) - \ln \det(\Omega) + \sum_{ij} \lambda'_{ij} |\Omega_{ij}| \quad (10)$$

$$s.t. \quad 0 \leq \mu < \lambda/m$$

Empirically we can solve (10) using graphical Lasso [16] with time complexity of $O(p^3)$, where $p = 21L$ is dimension of Ω . Nevertheless, if we want to generate m alternative solutions, the running time will be m times slower than PSICOV, i.e., generating a single solution.

Speedup. Here we describe a novel procedure to significantly reduce the computation time, by exploiting the facts that Ω is very sparse and there is still some similarity between two alternative solutions (although we want them to be diverse).

For a consecutive pair k and $k+1$, we can assume that the difference between λ^k and λ^{k+1} is not too large because of the update function for λ'_{ij} in (10). Therefore, we can estimate the new precision matrix Ω^{k+1} by first copying from Ω^k and then updating a small number of entries. Based on this observation, we propose following procedure.

In the first step, we apply the Quadratic Inverse Covariance (QUIC) method [17] to optimize (10) to generate the initial solution. Unlike the optimization methods used in

PSICOV, QUIC has a super linear convergence rate because it uses the second order information. At each iteration, QUIC approximates $F(\Omega)$ by its second-order Taylor expansion, computes the Newton direction, and descends at one coordinate.

In the second step, which is key for speedup, we identify which variables need to be updated. It is not hard to prove that if the following condition holds in current iteration, we do not need to update the submatrix Ω_{ij} in the remaining iterations.

$$|\nabla_{ij}F(\Omega)| \leq \lambda_{ij} \text{ and } \Omega_{ij} = 0 \quad (11)$$

Where $\nabla_{ij}F(\Omega)$ is the gradient of $F(\Omega)$ at position (i, j) . Let S_{fixed} denote the set of variables which do not need update and S_{free} the complement of S_{fixed} . Due to the similarity between two consecutive alternative solutions, we would expect that the number of variables in S_{free} is small. In practice we only need to update less than one percent of matrix entries at each step, which significantly improves the running time compared to simply running QUIC or graphical Lasso multiple times. That is, our method can quickly generate the alternative solutions after computing Ω^1 by the update rule in (11). We summarize our procedure in algorithm 1.

Algorithm 1 Algorithm for diverse M-Best contact maps

```

1: Input  $\lambda, \mu, m$ 
2: for  $k = 1, \dots, m$  do
3:   if  $k=1$  then
4:      $\Omega^1 = \text{QUIC}(S, \lambda)$ 
5:   else
6:     while not converge do
7:       Compute  $\lambda_{ij} = \lambda - \sum_{r=1}^{k-1} \mu \delta_0(\Omega^r)_{ij} \forall i, j$ 
8:       Identify  $S_{fixed}$  and  $S_{free}$  by (11)
9:       Compute the Newton direction  $\Delta$  on  $S_{free}$ 
10:      Compute step size  $\alpha$  by linear search
11:      Update  $\Omega^k = \Omega^k - \alpha \Delta$ 
12:    end while
13:  end if
14: end for
15: output  $\Omega^1, \Omega^2, \dots, \Omega^m$ 
```

To compute the update direction and step size in lines 9 and 10 of our algorithm, we need to maintain a temporary matrix to compute the Hessian matrix of $F(\Omega)$ efficiently. See [17] for more technical details. The algorithm runs very fast after computing Ω^1 due to step 8.

D. Contact Selection by Nuclear Norm

After computing the precision matrix, existing methods such as PSICOV select top predicted contacts based on the l_1 norm of each sub-matrix Ω_{ij} (see Eq.3), which can be interpreted as the interaction strength between two positions

(or columns). Here we propose a new method that uses the nuclear norm of each sub-matrix Ω_{ij} to measure the interaction strength of two positions. That is,

$$\tilde{\Omega}_{ij} = \sum_{i=1}^{20} \sigma_i \quad (12)$$

where σ_i is the i -th singular value of matrix Ω_{ij} . Again, we do not use the row and column in Ω_{ij} corresponding to the gap in calculating the nuclear norm. The nuclear norm is better than l_1 norm in capturing the sparsity pattern in each sub-matrix. Figure 1 shows two artificial 5×5 matrices with entry values ranging from 0 to 1 shown by different colors. The darker the color, the higher the value. Both matrices have the same l_1 norm. However, it is easy to see that the right matrix has a much stronger interaction pattern than the left one, which is consistent with their nuclear norms. The nuclear norms of the left and right matrices are 8.36 and 18.16, respectively. On average the nuclear norm is better than l_1 norm in measuring residue interaction strength since the former can better differentiate the sparsity pattern (right) from noise (left).

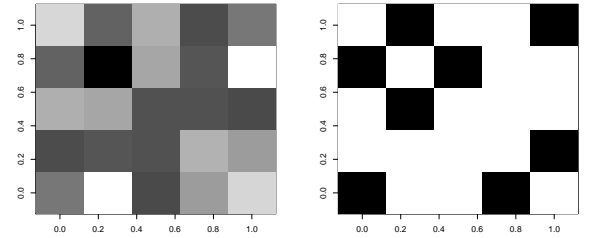


Figure 1: The figure shows two artificial matrices with different patterns, but same l_1 norm. The right one has a much stronger interaction pattern than the left one.

III. RESULTS

We evaluate our method by examining the improvement of the best solution over the first solution. We will also show that the best solutions generated by our method are also better than PSICOV and Evfold. We used three test sets: (1) The CASP10 test targets; (2) the PSICOV benchmark (150 Pfam families) and a subset of PDB25, which contain proteins with a large number of sequence homologs; and (3) a subset of PDB25 consisting of proteins with a small number of sequence homologs and thus, is very challenging. The PSICOV benchmark already contains multiple sequence alignment (MSA) from the Pfam database. For other test proteins, their multiple sequence alignments are generated using buildali.pl in the HHpred package [33], which runs PSI-BLAST [1] to find sequence homologs.

We consider only medium- and long-range contacts since they are much more important than short-range contacts.

Following PSICOV and other methods, we say there is a contact between two residues if their C_β Euclidean distance is less than 8\AA . A contact of two residues at sequence positions i and j is short-, medium- and long-range if $|i - j|$ is less than 12, between 12 and 24 or larger than 24, respectively.

Following PSICOV, we set the sparsity control parameter λ to 0.01. The diversity parameter μ is set to 0.0015 for all test cases. In this paper, for each protein under prediction, we predict only 5 different contact maps although it may lead to better accuracy by producing more solutions.

To ensure fair comparison, we employ the same pre- and post-processing procedures used by PSICOV, except that we use nuclear norm to measure the interaction strength rather than l_1 norm. Briefly, to reduce the impact of redundant sequences, we apply the same sequence weighting method as PSICOV. In particular, duplicate sequences are removed and columns containing more than 90% of gaps are deleted. The sequence is weighted using a threshold of 62% sequence identity. Similar to PSICOV and plmDCA [14], average-product correction (APC) [13] is applied to correct bias (see Eq. 4).

A. Performance on 123 CASP10 targets

1) Advantage of nuclear norm for contact selection:

We first show that nuclear norm is better in selecting top contacts. We run PSICOV to generate precision matrices and then employ nuclear norm and l_1 norm to select top $L/5$ contacts (L is sequence or MSA length), respectively. Both norms have similar performance in selecting medium-range contacts. The l_1 norm and the nuclear norm have accuracy 0.256 and 0.255, respectively. For long-range contacts, the nuclear norm and l_1 norm have accuracy 0.237 and 0.225, respectively.

2) *The best solution vs. the first solution:* When top $L/5$ predicted contacts are evaluated, the average medium-range prediction accuracy of the best (out of 5) solutions is 0.285, which is 11.51% better than the first solutions (0.255). The best solutions have long-range prediction accuracy 0.279, which is a 17.7% better than the first solutions (0.237). See Figure 2 for a detailed one to one comparison. When only the top $L/10$ predictions are evaluated, the first solutions have average medium- and long-range accuracy of 0.301 and 0.275, respectively. By contrast, the best solutions have average accuracy of 0.335 and 0.322, respectively. Overall, the best solutions tend to have a larger advantage over the first solutions in predicting long-range contacts, which is desirable since long-range contacts are the most useful for protein folding. In particular, for 32 out of the 123 test proteins, their best solutions have long-range accuracy at least 0.1 better than the first solutions. For 6 test proteins, their best solutions have long-range accuracy at least 0.2 better than the first ones. Note that the average accuracy of this test set is not very high mainly because the CASP10

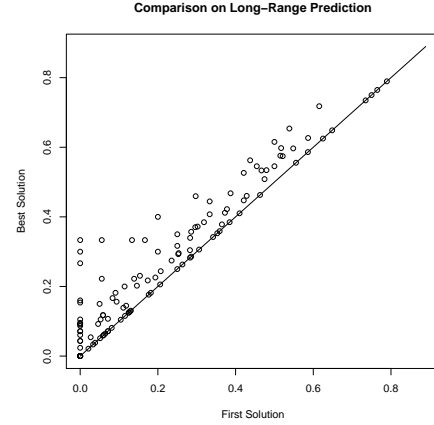


Figure 2: The figure illustrates the one-to-one comparison on top $L/5$ long-range contact prediction for CASP10 proteins. The X-axis is the accuracy for first solution and the Y-axis is that of the best solution.

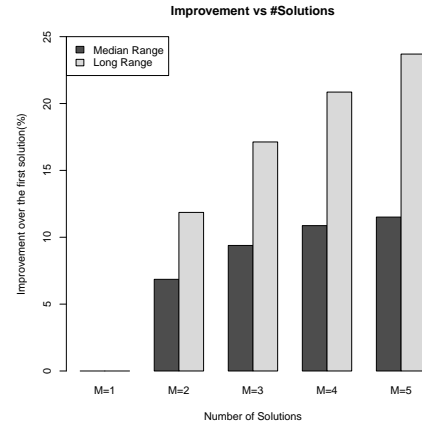


Figure 3: The accuracy improvement of the best solutions with respect to the number of solutions generated for a single protein. This figure shows that long-range prediction accuracy increases significantly with respect to the number of solutions.

set contains many test proteins with a very small number of sequence homologs.

3) *Performance with respect to the number of solutions:* As shown in Figure 3, the more solutions we generate, the better accuracy can be obtained by the best solutions. The long-range prediction accuracy can be improved by 24% when only 5 diverse solutions are generated. Even if only two solutions are generated for a protein, we can obtain 6.85% and 11.86% improvement, respectively, in medium- and long-range contact prediction.

4) *Performance with respect to the amount of homologous information:* Let M_{eff} denote the number of non-redundant

sequence homologs available for a protein (family) under prediction. Following [26], M_{eff} is calculated by

$$M_{eff} = \sum_i \frac{1}{\sum_j M_{i,j}} \quad (13)$$

where i and j represent two sequence homologs of the protein (family) under prediction and $M_{i,j}$ is a binary value indicating whether i and j are redundant. In particular, $M_{i,j} = 1$ if and only if i and j share more than 70% of sequence identity; otherwise 0. M_{eff} measures the amount of homologous information available for a protein (family), so it can be used to evaluate the difficulty level of a test protein (family).

We divide the 123 CASP10 test proteins into 10 groups by their $\ln M_{eff}$ values: $1 \sim 2, \dots, 9 \sim 10, \geq 10$, and then compute the average improvement in each group. As shown in Figure 4, the best out of 5 solutions have much larger advantage over the first solutions for proteins with smaller M_{eff} values. The reason is that with a limited amount of homologous information, the first predicted contact map may differ a lot from the ground truth.

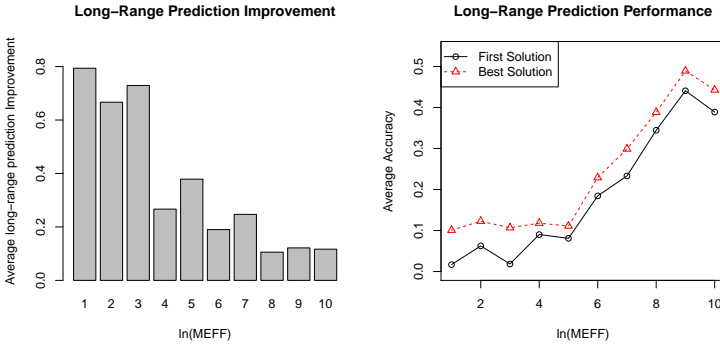


Figure 4: The left picture shows the relative improvement in long-range prediction accuracy of the best solutions over the first ones with respect to $\ln M_{eff}$, tested on the CASP10 set. The right picture shows the long-range accuracy of the best and first solutions.

B. Performance on proteins with a large number of sequence homologs

We build an easy data set, which consists of 149 test proteins in the PSICOV benchmark with $\ln M_{eff}$ between 7 and 10 and 140 proteins in PDB25 with $\ln M_{eff}$ between 5 and 7. The PDB25 list is downloaded from PISCES [37]. Any two proteins in PDB25 share less than 25% sequence identity. Compared to the CASP10 set, the test proteins in this data set on average have many more sequence homologs, so their contacts are relatively easier to predict.

Overall, when top $L/5$ predictions are evaluated, the best solutions are better than the first solutions by around 20% for both medium- and long-range prediction. More specifically,

for medium-range predictions, the best and first solutions have accuracy 0.381 and 0.315, respectively. For long-range prediction (see Figure 5), the best and first solutions have accuracy 0.464 and 0.391, respectively. In particular, for 111 out of 289 test proteins, the best solutions have accuracy at least 0.1 better than the first solutions for long-range prediction. For 23 test proteins, the best solutions have accuracy at least 0.2 better than the first ones for long-range prediction.

When top $L/10$ predictions are evaluated, the average accuracy of the first solutions are 0.380 and 0.419, respectively, for medium- and long-range contacts, while the best solutions have accuracy 0.445 and 0.499, respectively.

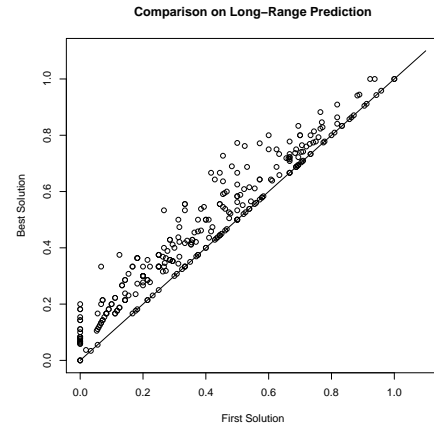


Figure 5: The above figure shows one-to-one comparison of the best and first solutions in terms of top $L/5$ long-range contact prediction accuracy for proteins with $\ln M_{eff}$ between 5 and 10.

C. Performance on proteins with a small number of sequence homologs

We build a data set of 122 test proteins from PDB25, all of which have $\ln M_{eff}$ smaller than 5. This set is much more challenging to predict since the test proteins have only limited amount of homologous information. When top $L/10$ predicted contacts are evaluated, the first solutions has average accuracy of 0.111 and 0.133 on medium-range and long-range contacts, respectively, while the best solutions have accuracy of 0.172 and 0.214, respectively. As shown in Figure 6, for 58 out of 122 test proteins, the best solutions have long-range accuracy at least 0.1 better than the first solutions. For 15 test proteins, the best solutions have long-range accuracy at least 0.2 better than the first ones.

When top $L/5$ predictions are evaluated, the average accuracy of the first solutions are 0.089 and 0.090, respectively, for medium- and long-range contacts, while the best solutions have accuracy 0.143 and 0.180, respectively.

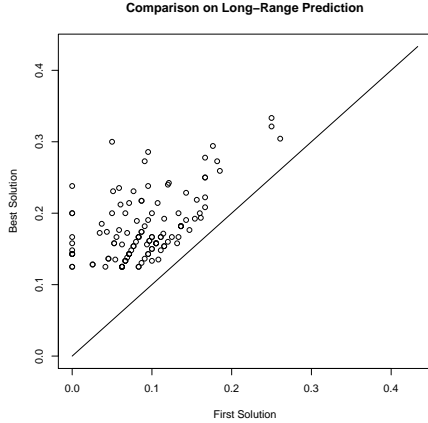


Figure 6: The above figure shows one-to-one comparison on top $L/10$ long-range contact prediction for proteins with $\ln M_{eff}$ less than 5.

Table I: Contact prediction accuracy of PSICOV, Evfold and our method on the CASP10 set.

	PSICOV	Evfold	Our M-best method
L/5 (long)	0.225	0.225	0.281
L/10 (long)	0.276	0.257	0.322
L/5 (median)	0.259	0.249	0.285
L/10 (median)	0.310	0.294	0.335

D. Comparison with other co-evolution-based methods

We compare the accuracy of our M-best method with a couple of popular co-evolution methods such as PSICOV and Evfold [26]. PSICOV is based upon Gaussian Graphical Model (GGM) and Evfold is derived by maximum entropy. We ran both PSICOV and Evfold by the default parameters suggested by their respective publications. Table I shows the average long-range prediction accuracy of PSICOV, Evfold and our method on the CASP10 set. As shown in this table, our M-best method outperforms PSICOV and Evfold. Our method is also based upon GGM, so its comparison with PSICOV implies that our method indeed can yield better accuracy by producing only a very small number of alternative solutions for an individual protein. In principle we may apply our M-best idea to the pseudo-likelihood model, which shall result in better accuracy than its corresponding tools such as plmDCA [14] and CCMpred [31] and further advance protein contact prediction.

E. Case Study

This section describes some specific cases on which the best solution is much better than the first one. One observation is that the best solutions tend to correctly predict more contacts in beta-sheet proteins. A specific example is PF07686 [4], which belongs to Immunoglobulin V-set

domain [30] and Ig-like domains resembling the antibody variable domain. V-set domains are found in diverse protein families, including immunoglobulin light and heavy chains. In addition, this domain belongs to SCOP family b.1.1.1 [2], which is an all-beta protein in greek-key shape that forms a sandwich consisting 7 strands in 2 sheets. The best out of 5 solutions correctly predicts 17 out of 22 long-range contacts, while the first solution only 11 of them.

Another example is PF01300, which belongs to Telomere recombination domain [36] and has been shown to bind preferentially to dsRNA. It has also been shown to be required for telomere recombination in yeast. This domain belongs to SCOP family d.115.1.1, which is an alpha and beta protein (a+b) in YrdC-like shape containing two additional strands in the C-terminal extension. The first solution correctly identifies 20 out of 35 long-range contacts, while our best solution identifies 27.

The CASP10 target T0688-D1 is a LRR (Leucine-rich repeat) domain [22] and a protein structural motif that forms an / horseshoe fold. It is composed of repeating 20-30 amino acid stretches that are unusually rich in the hydrophobic amino acid Leucine. In addition, this domain belongs to SCOP family c.10.2.1, which is capped at the N-end with a truncated EF-hand sub domain [27]. The best out of 5 solutions can correctly identify 17 out of 37 native long-range contacts while the first solution only 11 contacts.

IV. CONCLUSION & DISCUSSION

In this paper, we have presented a general structural learning method that can predict a set of diverse contact maps from an MSA. Compared to PSICOV, we only introduce one extra hyper-parameter to control the diversity of the predicted contact maps. Experiments confirm that our method can yield much better contact maps by generating only 5 different solutions, especially when the protein under consideration has a small number of sequence homologs. Our method works particularly well for contacts in beta-sheet proteins.

Currently we just fix the value of the hyper-parameter for the diversity control. In the future we are going to choose such a hyper-parameter depending on the number of non-redundant sequence homologs in a protein family. That is, when a family contains fewer sequence homologs, we shall allow more diversity due to less information available in the multiple sequence alignment. Otherwise, we shall allow less diversity.

Although our method generates only 5 different solutions, the best solutions usually have better accuracy than the first ones. One possible reason may be the fact that PSICOV model is biased since it assumes Gaussian distribution of an MSA, which is not true when only a small number of sequence homologs are available. It was noted that the solutions to the single MAP or maximize-likelihood problem is often biased, so providing multiple diverse solutions can

often help. Such phenomena occurred in a variety of computational biology problems, ranging from SNP data analysis [32], to sequence alignment [12], to protein threading [24], to protein structure alignment [38], [42], and to ab initio folding [15], [29]. The strategy we present here could be directly used to some of the above-mentioned problems and provide multiple informative solutions.

Note that in this paper our intention is not to develop the best contact prediction method that outperforms all the existing methods. Instead, this paper studies how we can make the best use of co-evolutionary information and advance the Gaussian Graphical model for contact prediction. To achieve the best contact prediction accuracy, it usually needs to combine both co-evolutionary and non-co-evolutionary information. For example, MetaPSICOV [19], the best contact predictor in CASP11 (the 11th Critical Assessment of Structure Prediction), is a supervised learning method that integrates four pure co-evolution-based methods (including PSICOV and Evfold) and a variety of non-co-evolutionary information. Such a supervised learning method may not work well on membrane proteins since it is trained by globular proteins. By contrast, a pure co-evolution-based unsupervised method can work on any proteins. Similar to the spirit of MetaPSICOV, we may also develop a supervised learning method to integrate our M-best method, other co-evolution methods and some non-co-evolutionary information to further advance contact prediction.

ACKNOWLEDGMENT

The authors are grateful to the financial support from National Institutes of Health R01GM0897532 and National Science Foundation CAREER award CCF-1149811

REFERENCES

- [1] Stephen F Altschul, Thomas L Madden, Alejandro A Schäffer, Jinghui Zhang, Zheng Zhang, Webb Miller, and David J Lipman. Gapped blast and psi-blast: a new generation of protein database search programs. *Nucleic acids research*, 25(17):3389–3402, 1997.
- [2] Antonina Andreeva, Dave Howorth, John-Marc Chandonia, Steven E Brenner, Tim JP Hubbard, Cyrus Chothia, and Alexey G Murzin. Data growth and its impact on the scop database: new developments. *Nucleic acids research*, 36(suppl 1):D419–D425, 2008.
- [3] Sivaraman Balakrishnan, Hetunandan Kamisetty, Jaime G Carbonell, Su-In Lee, and Christopher James Langmead. Learning generative models for protein fold families. *Proteins: Structure, Function, and Bioinformatics*, 79(4):1061–1078, 2011.
- [4] Alex Bateman, Lachlan Coin, Richard Durbin, Robert D Finn, Volker Hollich, Sam Griffiths-Jones, Ajay Khanna, Mhairi Marshall, Simon Moxon, Erik LL Sonnhammer, et al. The pfam protein families database. *Nucleic acids research*, 32(suppl 1):D138–D141, 2004.
- [5] Dhruv Batra, Payman Yadollahpour, Abner Guzman-Rivera, and Gregory Shakhnarovich. Diverse m-best solutions in markov random fields. In *Computer Vision–ECCV 2012*, pages 1–16. Springer, 2012.
- [6] Lukas Burger and Erik van Nimwegen. Disentangling direct from indirect co-evolution of residues in protein alignments. *PLoS computational biology*, 6(1):e1000633, 2010.
- [7] Jianlin Cheng and Pierre Baldi. Improved residue contact prediction using support vector machines and a large feature set. *BMC bioinformatics*, 8(1):113, 2007.
- [8] Simona Cocco, Remi Monasson, and Martin Weigt. From principal component to direct coupling analysis of coevolution in proteins: Low-eigenvalue modes are needed for structure prediction. *PLoS computational biology*, 9(8):e1003176, 2013.
- [9] David de Juan, Florencio Pazos, and Alfonso Valencia. Emerging methods in protein co-evolution. *Nature Reviews Genetics*, 14(4):249–261, 2013.
- [10] Pietro Di Lena, Piero Fariselli, Luciano Margara, Marco Vasura, and Rita Casadio. Is there an optimal substitution matrix for contact prediction with correlated mutations? *Computational Biology and Bioinformatics, IEEE/ACM Transactions on*, 8(4):1017–1028, 2011.
- [11] Pietro Di Lena, Ken Nagata, and Pierre Baldi. Deep architectures for protein contact map prediction. *Bioinformatics*, 28(19):2449–2457, 2012.
- [12] Chuong B Do, Samuel S Gross, and Serafim Batzoglou. ConTAlign: discriminative training for protein sequence alignment. In *Research in Computational Molecular Biology*, pages 160–174. Springer, 2006.
- [13] Stanley D Dunn, Lindi M Wahl, and Gregory B Gloor. Mutual information without the influence of phylogeny or entropy dramatically improves residue contact prediction. *Bioinformatics*, 24(3):333–340, 2008.
- [14] Magnus Ekeberg, Cecilia Lövkvist, Yueheng Lan, Martin Weigt, and Erik Aurell. Improved contact prediction in proteins: using pseudolikelihoods to infer potts models. *Physical Review E*, 87(1):012707, 2013.
- [15] Narayanan Eswar, Ben Webb, Marc A Marti-Renom, MS Madhusudhan, David Eramian, Min-yi Shen, Ursula Pieper, and Andrej Sali. Comparative protein structure modeling using modeller. *Current protocols in bioinformatics*, pages 5–6, 2006.
- [16] Jerome Friedman, Trevor Hastie, and Robert Tibshirani. Sparse inverse covariance estimation with the graphical lasso. *Biostatistics*, 9(3):432–441, 2008.
- [17] Cho-Jui Hsieh, Matyas A Sustik, Inderjit S Dhillon, and Pradeep D Ravikumar. Sparse inverse covariance matrix estimation using quadratic approximation. In *NIPS*, pages 2330–2338, 2011.

- [18] David T Jones, Daniel WA Buchan, Domenico Cozzetto, and Massimiliano Pontil. Psicov: precise structural contact prediction using sparse inverse covariance estimation on large multiple sequence alignments. *Bioinformatics*, 28(2):184–190, 2012.
- [19] David T Jones, Tanya Singh, Tomasz Kosciolk, and Stuart Tetchner. Metapsicov: combining coevolution methods for accurate prediction of contacts and long range hydrogen bonding in proteins. *Bioinformatics*, 31(7):999–1006, 2015.
- [20] David E Kim, Frank DiMaio, Ray Yu-Ruei Wang, Yifan Song, and David Baker. One contact for every twelve residues allows robust and accurate topology-level protein structure modeling. *Proteins: Structure, Function, and Bioinformatics*, 2013.
- [21] JL Klepeis and CA Floudas. Astro-fold: a combinatorial and global optimization framework for ab initio prediction of three-dimensional structures of proteins from the amino acid sequence. *Biophysical Journal*, 85(4):2119–2146, 2003.
- [22] Bostjan Kobe and Andrey V Kajava. The leucine-rich repeat as a protein recognition motif. *Current opinion in structural biology*, 11(6):725–732, 2001.
- [23] Alan Lapedes, Bertrand Giraud, and Christopher Jarzynski. Using sequence alignments to predict protein structure and stability with high accuracy. *arXiv preprint arXiv:1207.2484*, 2012.
- [24] Jianzhu Ma, Jian Peng, Sheng Wang, and Jinbo Xu. A conditional neural fields model for protein threading. *Bioinformatics*, 28(12):i59–i66, 2012.
- [25] Debora S Marks, Lucy J Colwell, Robert Sheridan, Thomas A Hopf, Andrea Pagnani, Riccardo Zecchina, and Chris Sander. Protein 3d structure computed from evolutionary sequence variation. *PloS one*, 6(12):e28766, 2011.
- [26] Faruck Morcos, Andrea Pagnani, Bryan Lunt, Arianna Bertolino, Debora S Marks, Chris Sander, Riccardo Zecchina, José N Onuchic, Terence Hwa, and Martin Weigt. Direct-coupling analysis of residue coevolution captures native contacts across many protein families. *Proceedings of the National Academy of Sciences*, 108(49):E1293–E1301, 2011.
- [27] Susumu Nakayama and Robert H Kretsinger. Evolution of the ef-hand family of proteins. *Annual review of biophysics and biomolecular structure*, 23(1):473–507, 1994.
- [28] Timothy Nugent and David T Jones. Accurate de novo structure prediction of large transmembrane protein domains using fragment-assembly and correlated mutation analysis. *Proceedings of the National Academy of Sciences*, 109(24):E1540–E1547, 2012.
- [29] Carol A Rohl, Charlie EM Strauss, Kira MS Misura, and David Baker. Protein structure prediction using rosetta. *Methods in enzymology*, 383:66–93, 2004.
- [30] Yoshinori Satow, Gerson H Cohen, Eduardo A Padlan, and David R Davies. Phosphocholine binding immunoglobulin fab mcpc603: An x-ray diffraction study at 2.7 Å. *Journal of molecular biology*, 190(4):593–604, 1986.
- [31] Stefan Seemayer, Markus Gruber, and Johannes Söding. Ccm-predfast and precise prediction of protein residue–residue contacts from correlated mutations. *Bioinformatics*, 30(21):3128–3130, 2014.
- [32] Mark Silberstein, Omer Weissbrod, Lars Otten, Anna Tzemach, Andrei Anisena, Oren Shtark, Dvir Tuberg, Eddie Galfrin, Irena Gannon, Adel Shalata, et al. A system for exact and approximate genetic linkage analysis of snp data in large pedigrees. *Bioinformatics*, 29(2):197–205, 2013.
- [33] Johannes Söding, Andreas Biegert, and Andrei N Lupas. The hhpred interactive server for protein homology detection and structure prediction. *Nucleic acids research*, 33(suppl 2):W244–W248, 2005.
- [34] Joanna I Sulkowska, Faruck Morcos, Martin Weigt, Terence Hwa, and José N Onuchic. Genomics-aided structure prediction. *Proceedings of the National Academy of Sciences*, 109(26):10340–10345, 2012.
- [35] Allison N Tegge, Zheng Wang, Jesse Eickholt, and Jianlin Cheng. Nncon: improved protein contact map prediction using 2d-recursive neural networks. *Nucleic acids research*, 37(suppl 2):W515–W518, 2009.
- [36] Marianna Teplova, Valentina Tereshko, Martin Egli, Ruslan Sanishvili, Andrzej Joachimiak, Tatyana Bushueva, and Wayne F Anderson. The structure of the yrdc gene product from escherichia coli reveals a new fold and suggests a role in rna binding. *Protein Science*, 9(12):2557–2566, 2000.
- [37] Guoli Wang and Roland L Dunbrack. Pisces: a protein sequence culling server. *Bioinformatics*, 19(12):1589–1591, 2003.
- [38] Sheng Wang, Jianzhu Ma, Jian Peng, and Jinbo Xu. Protein structure alignment beyond spatial proximity. *Scientific reports*, 3, 2013.
- [39] Zhiyong Wang and Jinbo Xu. Predicting protein contact map using evolutionary and physical constraints by integer programming. *Bioinformatics*, 29(13):i266–i273, 2013.
- [40] Martin Weigt, Robert A White, Hendrik Szurmant, James A Hoch, and Terence Hwa. Identification of direct residue contacts in protein–protein interaction by message passing. *Proceedings of the National Academy of Sciences*, 106(1):67–72, 2009.
- [41] Sitao Wu and Yang Zhang. A comprehensive assessment of sequence-based and template-based methods for protein contact prediction. *Bioinformatics*, 24(7):924–931, 2008.
- [42] Yang Zhang and Jeffrey Skolnick. Tm-align: a protein structure alignment algorithm based on the tm-score. *Nucleic acids research*, 33(7):2302–2309, 2005.